Advantages				
The method is efficient and easy to implement				
The resulting classifiers are easy interpretable. This does not happen for other approaches such as e.g. NNs				
Drawbacks				
The resulting classifier are seldom very effective [Cohen&Singer96, Joachims98,Lewis+96,Yang99]				
 If the documents in c_j occur in disjoint clusters a Rocchio classifier may be very ineffective: as all linear classifiers, it partition the space of documents in two linearly separable subspaces. But many naturally occurring categories are not linearly separable 				





Dip. di MatematicaF. Aiolli - Sistemi Informativi75Pura ed Applicata2006/2007

Instead of considering the set of negative training instances in its entirely, a set of near-positives might be selected (as in RF). This is called the query zoning method	
Near positives are more significant, since they are the most difficult to tell apart from the positives. They may be identified by issuing a Rocchio query consisting of the centroid of the positive training examples against a document base consisting of the negative training examples. The top-ranked ones can be used as near positives.	
Some claim that, by using query zones plus other enhancements, the Rocchio method can achieve levels of effectiveness comparable to state-of-the art methods while being quicker to train	

Linear	Classifiers)
--------	-------------	---

- A linear classifier is a classifier such that classification is performed by a dot product between the two vectors representing the document and the category, respectively. Therefore it consists in a document-like representation of the category c_i
- Linear classifiers are thus very efficient at classification time
- Methods for learning linear classifiers can be partitioned in two broad classes
 - Incremental methods (IMs) (or on-line) build a classifier soon after examining the first document, as incrementally refine it as they examine new ones
 - Batch methods (BMs) build a classifier by analyzing Tr all at once.









	Documents are	zero alono	i almost	all axes

- Most document pairs are very far apart (i.e., not strictly orthogonal, but only share very common words and a few scattered others)
- In classification terms: virtually all document sets are separable, for most ' any classification
- This is part of why linear classifiers are quite successful in this domain



kNN vs. Linear Classifiers

- Bias/Variance tradeoff
 - Variance ≈ Capacity
- kNN has high variance and low bias.
 - Infinite memory
- LCs has low variance and high bias.
 - Decision surface has to be linear (hyperplane)
- Consider: Is an object a tree?
 - Too much capacity/variance, low bias
 - Botanist who memorizes
 - Will always say "no" to new object (e.g., # leaves)
 - Not enough capacity/variance, high bias
 - Lazy botanist
 - Says "yes" if the object is green
 - Want the middle ground



